# IP Layer Metrics for 5G Edge Computing Service
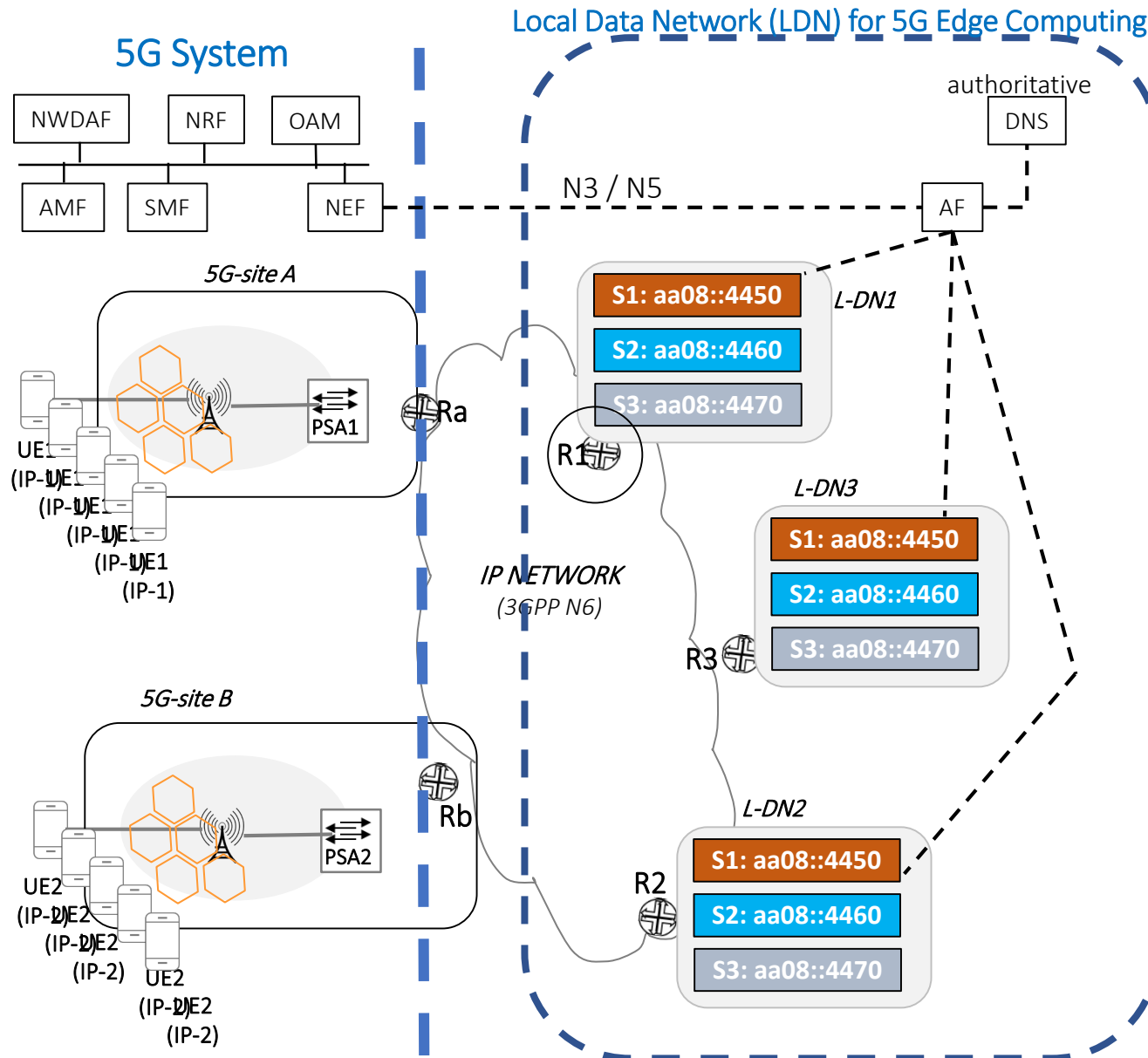
## draft-dunbar-ippm-5g-edge-compute-ip-layer-metrics-01

Linda Dunbar

HaoYu Song

John Kaippallimalil
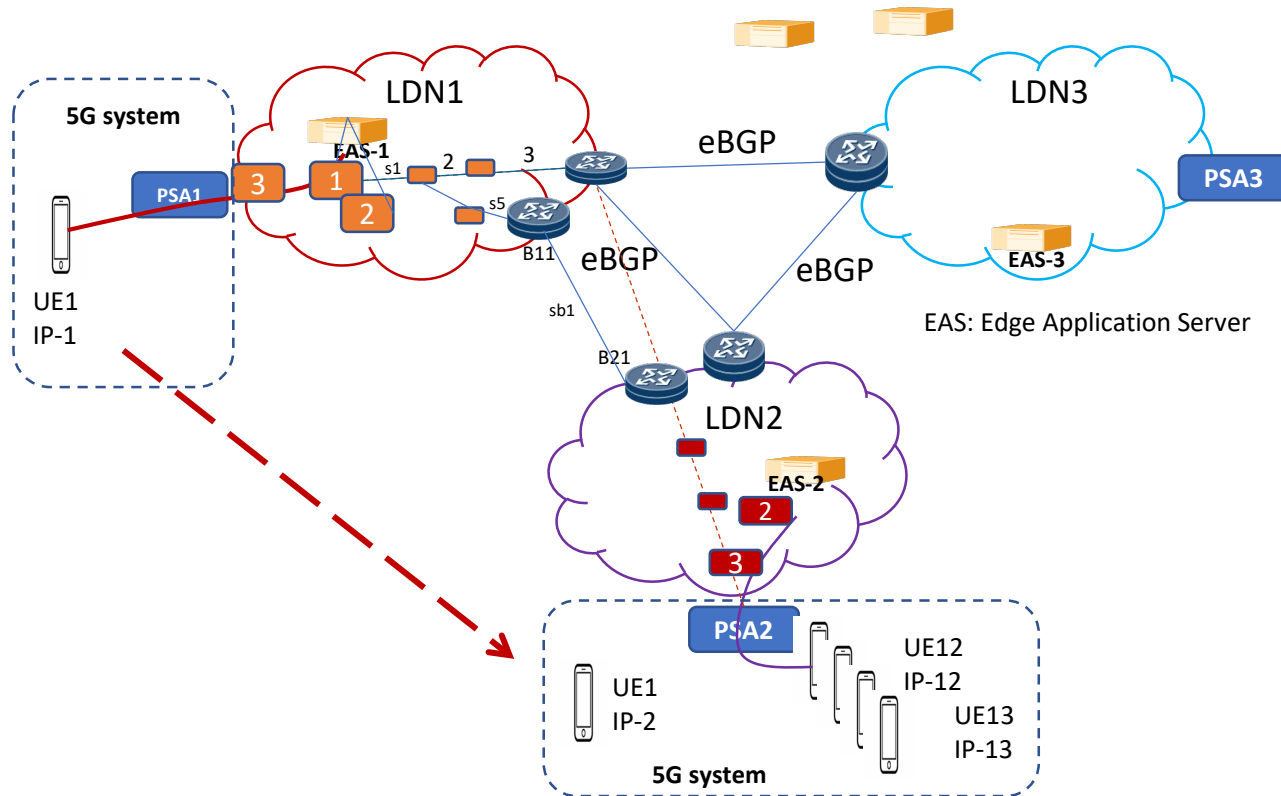
# 5G Edge Computing (3GPP TR23.748)

**Local Data Network (LDN) for 5G Edge Computing**

**5G System**

authoritative DNS

NWDAF  NRF  OAM

AMF  SMF  NEF

N3 / N5

AF

One Application has multiple Application Servers located in Edge Computing DCs

**5G-site A**

PSA1

Ra

UE1
(IP-1)
UE1
(IP-1)
UE1
(IP-1)
UE1
(IP-1)

R1

**L-DN1**

S1: aa08::4450
S2: aa08::4460
S3: aa08::4470

**L-DN3**

S1: aa08::4450
S2: aa08::4460
S3: aa08::4470

R3

*IP NETWORK*
*(3GPP N6)*

**5G-site B**

PSA2

Rb

UE2
(IP-2)
UE2
(IP-2)
UE2
(IP-2)
UE2
(IP-2)

R2

**L-DN2**

S1: aa08::4450
S2: aa08::4460
S3: aa08::4470

# From IP Network Perspective...

**ANYCAST: IP Layer Application ID -> multiple App servers**

**Benefit of using ANYCAST:**

- ✓ **dynamically load balance across locations based on network conditions.**
- ✓ **leverages the proximity information present in the network (routing) layer and**
- ✓ **eliminates the single point of failure and bottleneck at the DNS resolvers and application layer load balancers.**
- ✓ **removes the dependency on UEs using their cached destination IP addresses for extended period**



5G system

PSA1

UE1
IP-1

LDN1

EAS-1

s1   2   3

s5

B11

eBGP

sb1

B21

LDN2

EAS-2

2

3

PSA2

UE1
IP-2

UE12
IP-12

UE13
IP-13

5G system

eBGP

LDN3

PSA3
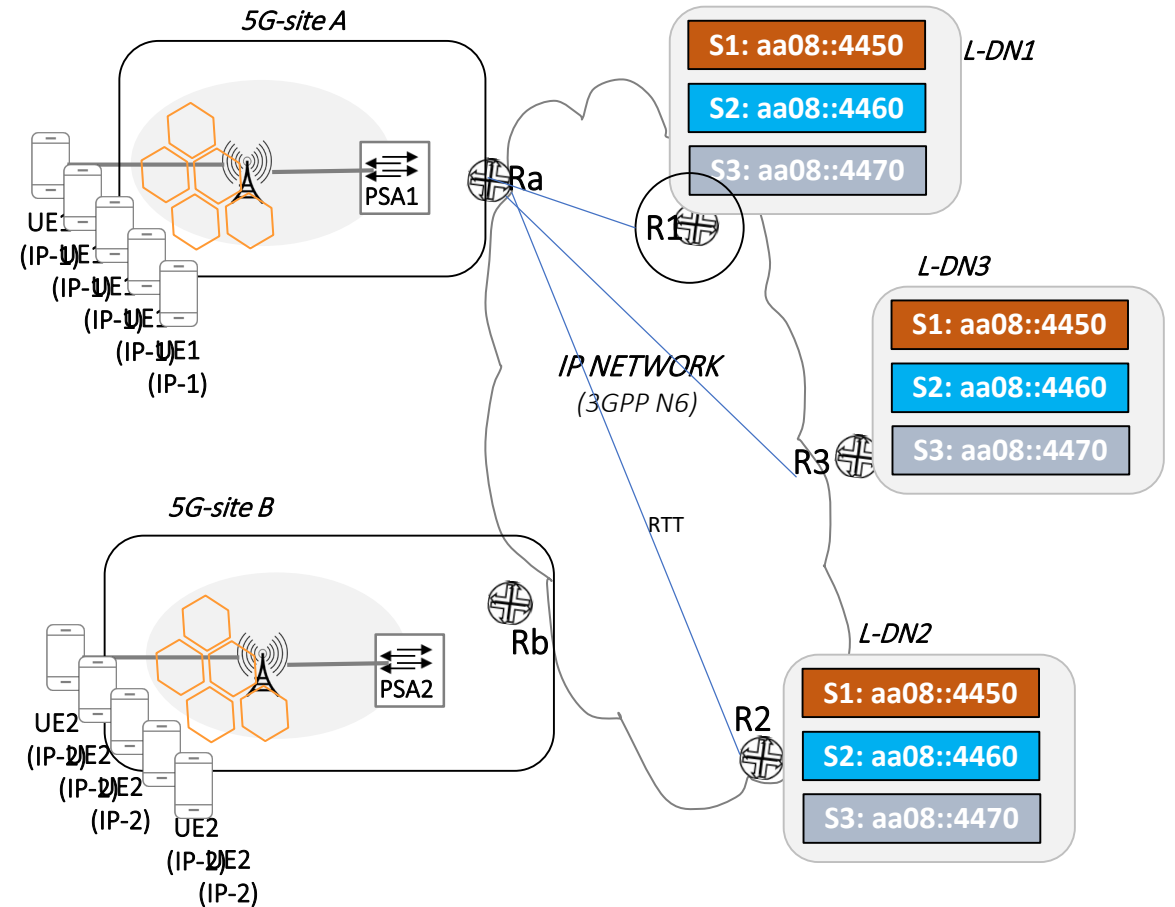
EAS-3

eBGP

eBGP

EAS: Edge Application Server

Problem 1: Selecting 5G Edge Application Location

- ➤ Many mini data centers can be close in proximity, making it difficult to differentiate in Routing Hops for App servers hosted in them,
- ➤ Some data centers can have higher capacity than others,
- ➤ Some sites may be more preferred when a UE anchored to a new 5G Site

Problem 2: UE mobility creates unbalanced anycast distribution

# RTT to an ANYCAST Address in 5G EC

- RTT to "app.net" ANYCAST S1:
- List of {
  - R1: RTT value
  - R2: RTT value
  - R3: RTT value
  }

# Algorithm in Selecting the Optimal Target Location

To compare the cost to reach the Application Server between the Site-i or Site-j:

$$\text{Cost-}i = \min\left(w * \left(\frac{\text{Load-}i * \text{CP-}j}{\text{Load-}j * \text{CP-}i}\right) + (1-w) * \left(\frac{\text{Pref-}j * \text{Delay-}i}{\text{Pref-}i * \text{Delay-}j}\right)\right)$$

- Load-i: Load Index at Site-I = w1*ToPackets+w2*FromPackes+w3*ToBytes+w4*FromBytes

  0<= wi <=1 and w1+ w2+ w3+ w4 = 1.

- CP-i (Capacity-i) (higher value means higher capacity): capacity index at the site i.

- Delay-i: Network latency measurement (RTT) to the A-ER that has the Application Server attached at the site-i.

- Pref-i (Preference Index: higher value means higher preference): Network Preference index for the site-I.

- w: Weight for load and site information,

  - 0<= w <=1: If smaller than 0.5, Network latency and the site Preference have more influence; otherwise, Server load and its capacity have more influence.

# Next step:

- To standardize the IP Layer Metrics for Application Servers
    - To make network more aware of application server running status and environment
    - To achieve more optimized delivery of service
- Need your feedback